

A General and Parallel Platform for Mining Co-Movement Patterns over Large-scale Trajectories

ABSTRACT

1. INTRODUCTION

The prevalence of positioning devices has drastically boosted the scale and spectrum of trajectory collection to an unprecedented level. Tremendous amounts of trajectories, in the form of sequenced spatial-temporal records, are continually generated from animal telemetry chips, vehicle GPSs and wearable devices. Data analysis on large-scale trajectories benefits a wide range of applications and services, including traffic planning [1], animal analysis [2], and social recommendations [3], to name just a few.

A crucial task of data analysis on top of trajectories is to discover co-moving patterns. A *co-movement* pattern [4] refers to a group of objects traveling together for a certain period of time and the group is normally determined by spatial proximity. A pattern is prominent if the size of the group exceeds M and the length of the duration exceeds K , where M and K are parameters specified by users. Rooted from such basic definition and driven by different mining applications, there are a bunch of variants of co-movement patterns that have been developed with more advanced constraints.

Table 1 summarizes several popular co-moving patterns with different constraints in the attributes of clustering in spatial proximity, consecutiveness in temporal duration and computational complexity. In particular, the *flock* [5] and *group* [6] patterns require all the objects in a group to be enclosed by a disk with radius r ; whereas the *convoy* [7], *swarm* [8] and *platoon* [9] patterns resort to density-based spatial clustering. In the temporal dimension, *flock* [5] and *convoy* [7] require all the timestamps of each detected spatial group to be consecutive, which is referred to as *global consecutiveness*; whereas *swarm* [8] does not impose any restriction. The *group* [6] and *platoon* [9] adopt a compromised manner by allowing arbitrary gaps between the consecutive segments, which is called *local consecutiveness*. They introduce a parameter L to control the minimum length of each local consecutive segment.

	Proximity	Consecutiveness	Time Complexity
flock [10]	disk-based	global	$O(\mathbb{O} \mathbb{T} ^2 + \mathbb{O} \log(\mathbb{O}))$
convoy [7]	density-based	global	$O(\mathbb{O} ^2 + \mathbb{O} \mathbb{T})$
group [6]	disk-based	local	$O(\mathbb{O} ^2 \mathbb{T})$
platoon [9]	density-based	local	$O(2^{ \mathbb{O} } \mathbb{O} \mathbb{T})$
swarm [8]	density-based	no constraint	$O(2^{ \mathbb{O} } \mathbb{O} \mathbb{T})$

Table 1: Constraints and complexity of co-movement patterns. The time complexity indicates the performance in the worst case, where $|\mathbb{O}|$ is the total number of objects and $|\mathbb{T}|$ is the number of discretized timestamps.

To better illustrate the temporal constraints of these co-movement patterns, we design a scenario with six example trajectories in Figure 1. In the setting, the temporal dimension is discretized into six snapshots. In each snapshot, objects are clustered as shown in the dotted circles. We use M to denote the minimum size of a object group, K to denote the minimum number of timestamps in a duration, L to denote the minimum size of local-consecutive segments in a duration.

Global consecutiveness: *Flock* and *convoy* patterns enforce the pattern duration to be *global consecutive*. We discover one *flock* and one *convoy* pattern, which both correspond to $\{o_3, o_4 : 1, 2, 3\}$. It can be verified that no other object groups appears in more than three consecutive timestamps.

Local consecutiveness: *Group* and *platoon* patterns require the duration to be *locally consecutive*. When setting $L = 2$, we observe three *group* patterns and three *platoon* patterns. Specifically, *group* aims to find the pattern with maximal object set whose duration conforms to L consecutiveness. Taking the pattern $\{o_3, o_4, o_5 : 2, 3\}$ as an example. $\{o_3, o_4, o_5\}$ is indeed maximal. If another object, say o_2 , is added to the group, the duration reduces to $\{3\}$ which is not L -consecutive. In fact, adding any other object to this group invalidates the L -consecutiveness. In contrary, *platoon* finds the maximal object set whose duration is not only L -consecutive but also with size no less than K . Therefore, *platoon* finds a different set of patterns. For example, $\{o_4, o_5 : 2, 3, 5, 6\}$ is a *platoon* pattern but not a *group* pattern. This is because the object set $\{o_4, o_5\}$ is not maximal wrt. *group*, as we can see that adding o_3 to the object set results in a larger *group* pattern. Therefore, $\{o_4, o_5 : 2, 3, 5, 6\}$ is not a *group*.

No consecutiveness: *Swarm* pattern relaxes the consecutiveness in the duration. Thus as shown in the Figure, *swarm* discover the most number of patterns. In particular,

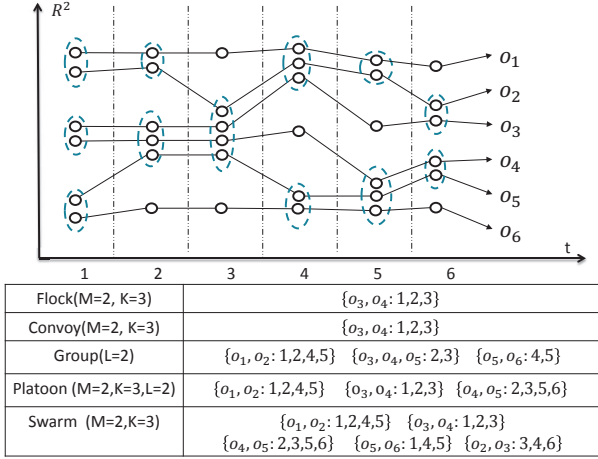


Figure 1: Trajectories and co-movement patterns; The example consists of six trajectories across six snapshots. The clusters of objects in each snapshots are denoted using dotted circles. Each type of pattern has its parameters. M denotes the minimum size of objects. K denotes the minimum size of timestamps. L denotes the minimum size of locally consecutive timestamps.

given the same M and K , *swarm* result covers all patterns discovered in *flock*, *convoy* and *platoon*. Furthermore, it finds more patterns that are not covered by *flock*, *convoy* and *platoon*. For instance, $\{o_5, o_6, 1, 4, 5\}$ is a *swarm* but is neither a *platoon* (with $L = 2$) nor a *convoy*.

As can be seen, there are various co-movement patterns requested by different applications and it is cumbersome to design a tailored solution for each type. We further identify that existing patterns either miss out interesting patterns (i.e., *convoy*) or likely contain non-interesting patterns (i.e., *swarm*, *platoon*). Hence, there is a particular demand for a new pattern model to fine control the pattern result without losing expressiveness. The other issue with existing methods is that they are built on top of centralized indexes that are not scalable. To the best of our knowledge, the maximum number of trajectories ever evaluated is up to hundreds of thousands of trajectories. In practice, it is rather common to collect at least millions of trajectories and their scalability is left unknown. We conduct a theoretical analysis on the worst-case complexity (as listed in Table 1 as well as an experimental evaluation with million-scale cardinality of a trajectory database (as shown in Figure ??). Results show that their performances degrade dramatically as the dataset size scales up and none of the existing solutions is scalable to handle large-scale trajectories.

Therefore, our primary contributions in this paper are to close these two gaps. First, we summarize two anomalies from existing patterns, namely, *missing-pattern* and *loose-connection*. The *missing-pattern* [8] anomaly arises due to the stringent constraints on the duration of a pattern. As shown in Figure 2 (a), if the duration size is set to be $K = 4$, neither *flocks* nor *convoys* can be discovered. As we notice, object o_1 is temporally far from o_2 at timestamp 4, which is likely to be the result of errors in interpretation of missing points, or o_1 faces traffic control at time 4. Such an anomaly can be resolved by *swarm* and *platoon* due to a relaxed constraint on the duration. However, *swarm* and

platoon's relaxations encompass a type of non-interesting patterns which is referred as *loose-connection* [9] anomaly. As shown in Figure 2 (b), the two objects o_1, o_2 would form a *platoon* pattern $\{o_1, o_2: 1, 2, 3, 102, 103, 104\}$. However, the consecutive segment on the pattern's duration is 98 timestamps apart, making the co-movement behavior very weak. In reality, such a pattern is likely to be induced by periodic movements of unrelated objects such as, vehicles stopping at the same petrol station, animals pausing at the same water source etc. In current literature, users are unable to explicitly exclude the loosely-connected patterns even when those patterns are unwanted. To cope with both of the two anomalies, we propose the *general co-movement pattern* by introducing the gap parameter G , which requires the distance of timestamps in the duration to be no larger than G . By so doing, we gain a fine-grained control over the pattern duration, which alleviates both the *missing-pattern* and the *loose connection* problem. Meanwhile, the general co-movement pattern does not concede its expressiveness. As we show in later sections, the general co-movement pattern is able to express existing patterns by setting appropriate parameters. Therefore, users are still able to exclude non-consecutive patterns or include loosely-connected patterns when they feel necessary.

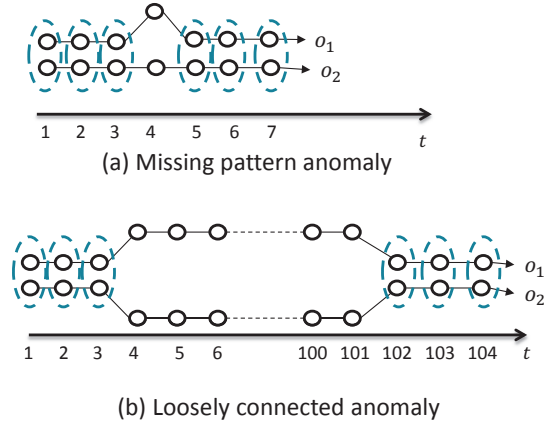


Figure 2: Two anomalies in existing patterns. (a) *Missing-pattern* anomaly in *flock* and *convoy*. When $K = 4$, none of the two patterns can be discovered. (b) *Loose-connection* anomaly in *platoon* and *swarm*. The consecutive segment of o_1 and o_2 are 98 timestamps apart, however, the pattern $\{o_1, o_2: 1, 2, 3, 102, 103, 104\}$ is included in *platoon* and *swarm* results.

Second, we propose a parallel solution on modern MapReduce platforms for scalable pattern mining. The major challenge in designing MapReduce-based algorithms is to make proper partitions of input data. In GCMP mining, we enforce both the *soundness* and the *completeness* of the partitions. Such properties ensure that neither false-patterns nor miss-patterns are possible in our solution. To meet such partition requirements as well as to keep the shuffle amount to a minimum, we first design a naive *Temporal Replication and Mining* (TRM) approach, which partitions trajectories into groups of consecutive snapshots. Then, we design a line-sweep method for mining GCMP from each partition. We prove that in TRM, the partition is complete and sound when a snapshot is replicated $O(|T|)$ times. Then, we design

a novel *Star Partition and Mining* (SRM) approach which significantly reduces the data shuffled as compare to TRM. In SRM, we design a conceptual connection graph based on proximity among objects. We adapt a *star partition* which cut the graph by replicating vertices. Afterwards, we design an Apriori-like method to mine GCMP in each partition. We prove the correctness of SPM and show that total data been replicated is $O(|\mathbb{O}|)$. Moreover, we utilize *temporal monotonicity* in GCMP to further reduce the shuffling and mining cost in SPM. Lastly, we adapt various engineering level techniques to support efficiently deploying our algorithms in Apache Spark which is one of the most popular MapReduce platforms.

We conduct a set of extensive experiments on XXX datasets with million-scale trajectories. The results show that XXX.

The rest of our paper is organized as follows: Section 2 summarizes the relevant literature on trajectory pattern mining; Section 3 forms the definition of the general co-movement pattern mining; Section 4 presents our parallel architecture; The solution of mining the general co-movement pattern mining is presented in Section 5 and Section 6. Section 7 discuss various optimization techniques to boost the system performance; Section 8 conducts extensive experiments to showcase the usefulness and efficiency of our system and finally Section 9 concludes our paper.

2. RELATED WORKS

The *co-movement patterns* in literature consist of five members, namely *group* [6], *flock* [10], *convoy* [7], *swarm* [8] and *platoon* [9]. We have demonstrated the semantics of these patterns in Table 1 and Figure 1. In this section, we focus on comparing the techniques used in these works. For more trajectory patterns other than *co-movement patterns*, interested readers may move to [11] for a comprehensive survey.

2.1 Flock and Convoy

The difference between *flock* and *convoy* lies in the object clustering methods. In *flock* objects are clustered based on their distance. Specifically, the objects in the same cluster needs to have a pair-wised distance less than *min_dist*. This essentially requires the objects to be within a disk-region of delimiter less than *min_dist*. In contrast, *convoy* cluster the objects using density-based clustering [12]. Technically, *flock* utilizes a m^{th} -order Voronoi diagram [13] to detect whether a subset of object with size greater than m stays in a disk-region. *Convoy* employs a trajectory simplification [14] technique to boost pairwise distance computations in the density-based clustering. After clustering, both *flock* and *convoy* use a line-sweep method to scan each snapshots. During the scan, the object group appears in consecutive timestamps is detected. Meanwhile, the object groups that do not match the consecutive constraint are pruned. However, such a method faces high complexity issues when supporting other patterns. For instance, in *swarm*, the candidate set during the line-sweep grows exponentially, and many candidates can only be pruned after the entire snapshots are scanned.

2.2 Group, Swarm and Platoon

Different from *flock* and *convoy*, all the *group*, *swarm* and *platoon* patterns have more constraints on the pattern duration. Therefore, their techniques of mining are of the same skeleton. The main idea of mining is to grow object set from

an empty set in a depth-first manner. During the growth, various pruning techniques are provided to prune unnecessary branches. *Group* pattern uses the Apriori property among patterns to facilitate the pruning. *Swarm* adapts two more pruning rules called backward pruning and forward pruning. *Platoon* further adapts a prefix table structure to guide the depth-first search. As shown by Li et.al. [9], *platoon* outperforms other two methods in efficiency. However, the three patterns are not able to directly discover the general co-movement pattern. Furthermore, their pruning rules heavily rely on the depth-first search nature, which lost its efficiency in the parallel scenario.

THESE WORKS ARE MOST RELATED TO OUR PROBLEMS, SO I REMOVED OTHER RELATED WORKS FOR NOW.

3. DEFINITIONS

Let $\mathbb{O} = \{o_1, o_2, \dots, o_n\}$ be the set of objects and $\mathbb{T} = \{1, 2, \dots, m\}$ be the descritized temporal dimension. A time sequence T is defined as a subset of \mathbb{T} , i.e., $T \subseteq \mathbb{T}$, and we use $|T|$ to denote sequence length. Let T_i be i -th entry in T and we say T is consecutive if $\forall 1 \leq i \leq |T|-1, T_{i+1} = T_i + 1$. It is obvious that any time sequence T can be decomposed into consecutive segments and we say T is L -consecutive [9] if the length of all the consecutive segments is no smaller than L .

As illustrate in Figure 2, patterns adapting the notion of L -consecutiveness (e.g., *platoon* and *group*) still suffer from *loose connection* problem. To avoid such an anomaly without losing pattern generality, we introduce a parameter G to control the gaps between timestamps in a pattern. Formally, a G -connected time sequence is defined as follows:

Definition 1 (G -connected). *A time sequence T is G -connected if the gap between any of its neighboring timestamps is no greater than G . That is $\forall T_i, T_{i+1} \in T, T_{i+1} - T_i \leq G$.*

We take $T = \{1, 2, 3, 5, 6\}$ as an example, which can be decomposed into two consecutive segments $\{1, 2, 3\}$ and $\{5, 6\}$. T is not 3-consecutive since the length $\{5, 6\}$ is 2. Thus, it is safe to say either T is 1-consecutive or 2-consecutive. On the other hand, T is 2-connected since the maximum gap between its neighboring time stamps is $5 - 3 = 2$. It is worth noting that T is not 1-connected because the gap between T_3 and T_4 is 2 (i.e., $5 - 3 = 2$).

Given a trajectory database descritized into snapshots, we can conduct a clustering method, either disk-based or density-based, to identify groups with spatial proximity. Let T be the set of timestamps in which a group of objects O are clustered. We are ready to define a more general co-movement pattern:

Definition 2 (General Co-Movement Pattern). *A general co-movement pattern finds a set of objects O satisfying the following five constraints: 1) closeness: the objects in O belong to the same cluster in the timestamps of T ; 2) significance: $|O| \geq M$; 3) duration: $|T| \geq K$; 4) consecutiveness: T is L -consecutive; and 5) connection: T is G -connected.*

There are four parameters in our general co-movement pattern, including object constraint M and temporal constraints K, L, G . By customizing these parameters, our pattern can express other patterns proposed in previous literature, as illustrated in Table 2. In particular, by setting $G = |T|$, we achieve the *platoon* pattern. By setting

Pattern	M	K	L	G	Clustering
Group	2	1	2	$ T $	Disk-based
Flock	\cdot	\cdot	K	1	Disk-based
Convoy	\cdot	\cdot	K	1	Density-based
Swarm	\cdot	\cdot	1	$ T $	Density-based
Platoon	\cdot	\cdot	\cdot	$ T $	Density-based

Table 2: Expressing other patterns using GCMP. \cdot indicate a user specified value. M represents the object size constraints. K represents duration constraint. L represents consecutiveness constraint. G represents the connection constraints.

$G = |T|, L = 1$, we achieve the *swarm* pattern. By setting $G = |T|, M = 2, K = 1$, we gain the *group* pattern. Finally by setting $G = 1$, we achieve the *convoy* and *flock* pattern. In addition to the flexibility of representing other existing patterns, our GCMP is able to avoid the *loose connection* anomaly by tuning the parameter G . It is notable that GCMP cannot be modeled by existing patterns. AS MENTIONED IN WECHAT, POLISH THIS PART.

It is also observable that the number of patterns in GCMP could be exponential under some parameter settings (i.e., when expressing *swarm*). In particular, given a parameter M , if a pattern P is valid, then any subset of P with size M is also a valid pattern. This results in additional $\sum_{M \geq i \geq |P.O|} \binom{|P.O|}{i}$ patterns, which is clearly overwhelming and redundant. For all these patterns, output P is sufficient. Therefore, we define the *Closed General Co-Movement Pattern* as follows:

Definition 3 (Closed General Co-Movement Pattern). A general co-moving pattern $P = \langle O : T \rangle$ is closed if and only if there does not exist another general co-moving pattern P' s.t. $P.O \subseteq P'.O$.

For example, let $M = 2, K = 2, L = 1, G = 1$. In Figure 1, the pattern $P_1 = \{o_3, o_4 : 1, 2, 3\}$ is not a closed pattern. This is because $P_2 = \{o_3, o_4, o_5 : 2, 3\}$ is a closed pattern since $P_2.O \supset P_1.O$. The closed pattern avoids outputting duplicate information, thus making the result patterns more compact.

Our definition of GCMP is free from clustering method. Users are able to supply different clustering methods to facilitate different application needs. We currently expose both disk-region based clustering and DBSCAN as options to the user.

In summary, the goal of this paper is to present a parallel solution for discovering closed GCMP from large-scale trajectory data.

Before we move on to the algorithmic part, we list the notations that are used in the following sections.

4. OVERVIEW OF MINING GCMP IN PARALLEL

We adapt the MapReduce paradigm for designing a parallel solution of mining GCMP. In this section, we briefly describe the preliminaries on MapReduce and then describe the overview of our framework in mining GCMP.

4.1 Preliminary on MapReduce

Symbols	Meanings
Tr_i	Trajectory of object i
S_t	Snapshot of objects at time t
\mathbb{O}	Set of objects
M	Object size constraint
K	Duration constraint
L	Consecutiveness constraint
G	Connection constraint
T	Time sequence
$C_t(o)$	the cluster of object o at time t
Sr_i	The star structure of object i

Table 3: Notions that will be used

MapReduce (MR) was formally proposed by Dean et. al. [15] and has subsequently implemented by many open source systems. Those systems provide handy APIs with fault tolerances and are popularly used as large-scale data processing platforms. In simple words, there are two conceptual types of computing nodes in MR, namely the *mappers* and the *reducers*. The execution of a MR algorithm consists of two stages: First, input data are partitioned and read by a *map* function on each mapper. Mappers emit key-value pairs which are *shuffled* over the network to reducers. Second, reducers process the received data using a *reduce* function then write the outputs.

Despite the simpleness of the paradigm, there are two concerns raised in designing MR algorithms. First, since reducers are required to be independent, partitioning data to fit the independence could be challenging. Second, since the *shuffle* stage requires network access, the data been shuffled should be minimized. We take these concerns in consideration when designing our solutions for mining GCMP.

4.2 MapReduce Processing for Mining GCMP

Our GCMP mining process consists of two MR jobs, which are further divided into four MR stages as illustrated in Figure 3. The first MR job is to cluster objects at each snapshot (i.e., $\forall t, o$, compute $C_t(o)$). As shown in Figures 3(a)-(b), in map phase, trajectories are first reorganized by timestamps and object locations at the same timestamps form a cluster. In reduce phase, clustering of objects at each snapshot is processed independently. The second MR job is to mine GCMPs from clusters in each snapshot. We design two MR algorithms (i.e., TRM and SPM) for the purpose of comparison. Both the two algorithms would partition snapshots in map phase (as in Figure 3(c)) and mine GCMPs from each partition in reduce phase (as in Figure 3(d)).

Although we need two MR jobs to complete the GCMP mining task, it is easy to pipeline the two jobs to exploit data locality. Specifically, the reducer output at step (b) can be directly reused as the input to the mappers at step (c). Therefore we do not need to transfer data between the two jobs. Modern MR platforms, especially Spark, have already supported such a kind of pipeline by default.

We note that the first MR job is easy to design since each reducer only needs one snapshot for clustering. In contrast, it is challenging to design the second job. This is because valid patterns may spray across multiple snapshots or contain different object sets, where inappropriate partitioning of snapshots may fail to discover certain valid patterns. Formally, a valid partition strategy needs to meet the following

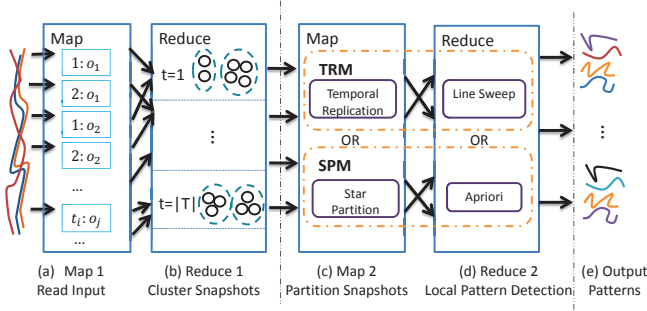


Figure 3: System flow of mining GCMP. (a)(b) correspond to the first MR job which compute the clusters at each snapshot; (c)(d) correspond to the second MR jobs which mines GCMP in parallel.

requirements: (a) the resulted partitions need to preserve enough information so that real patterns can be discovered in the reduce phase. (b) the resulted partitions need to ensure that the patterns discovered in the reduce phase are valid patterns so that no further verification is required. We formalize these two properties as *completeness* and *soundness* as follows:

Definition 4 (Completeness and Soundness). *Let a partition method \mathbb{P} partitions original trajectories Tr into multiple parts, Par_1, \dots, Par_m . \mathbb{P} is complete if for every pattern P that is valid in Tr , $\exists Par_i$ s.t. P is valid in Par_i . \mathbb{P} is sound if for all patterns that are valid in any Par_i , they are also valid in Tr .*

The completeness ensures that no true patterns are missed out. The soundness ensures that no false patterns are reported. If a partition method is both sound and complete, then it can be used in the second MR job to facilitate GCMP mining.

Apparently, replicating the entire trajectories to each partition meets the *soundness* and *completeness* requirements. However, it burdens the network shuffle and limits the parallelism. Our objective is thus to design a complete and sound partition method that minimize the network shuffles. In the following sections, we describe a naive *temporal-based* partition-and-mining method called *Temporal Replication and Mining* (TRM) towards a parallel solution of GCMP mining. Then, we present a novel *object-based* partition-and-mining method called *Star Partition and Mining* (SPM) which resolves the deficiencies of TRM method.

5. TEMPORAL REPLICATION AND MINING

The idea of *Temporal Replication and Mining* (TRM) is to group temporally closed snapshots together, such that patterns can be mined on each group of the snapshot. In order to achieve the *completeness* and *soundness* during partitioning, we allow replication of snapshots among groups. The TRM is described in Algorithm 1.

As shown in Algorithm 1, the TRM algorithm takes three steps. First, in the map phase, each snapshot is keyed with its timestamp (lines 1-6). Second, in the partition phase, every snapshot is grouped with its next $(\lfloor \frac{K}{L} \rfloor - 1) * G + K + L$ snapshots to form a partition (lines 7-11). We will

Algorithm 1 Temporal Replication and Mining

Require: list of $\langle t, S_t \rangle$ pairs

- 1: —Map Phase—
- 2: **for all** $\langle t, S_t \rangle$ **do**
- 3: **for all** $i \in 1 \dots (K - 1) * G + K$ **do**
- 4: emit a $\langle t - i, S_t \rangle$ pair
- 5: **end for**
- 6: **end for**
- 7: —Partition and Shuffle Phase—
- 8: **for all** $\langle t, S \rangle$ pair **do**
- 9: group-by t , emit a $\langle t, Par_t \rangle$,
- 10: where $Par_t = \{S_t, S_{t+1}, \dots, S_{t+(\lfloor \frac{K}{L} \rfloor - 1) * G + K + L}\}$
- 11: **end for**
- 12: —Reduce Phase—
- 13: **for all** $\langle t, Par_t \rangle$ **do**
- 14: lineSweepMining(Par_t)
- 15: **end for**

shortly discuss how the group size is determined. Third, in the reduce phase, a lineSweepMining method is invoked to mine GCMP within each partition (lines 12-14). It is easy to see that this method replicates a snapshots at most $(\lfloor \frac{K}{L} \rfloor - 1) * G + K + L$ times.

5.0.1 Temporal Replication Partition

The size of replication is critical for the performance of TRM algorithm. If the size of replication is too large, the shuffle cost as well as the reduce cost would be high. On the contrary, if the size of replication is too small, the *completeness* and *soundness* properties cannot be satisfied. In the Algorithm 1, the partition size is chosen as $(\lfloor \frac{K}{L} \rfloor - 1) * G + K + L$. As stated in the following theorem, such a partition method is sound and complete.

Theorem 1 (Soundness and Completeness of Replication). *Let \mathbb{P} be as follows: for each snapshot S_t , create a partition $Par_t = \{S_t, \dots, S_{t+(\lfloor \frac{K}{L} \rfloor - 1) * G + K + L}\}$. Then \mathbb{P} is sound and complete.*

Proof. The soundness of partition is trivial. Given a valid pattern P , let $T' \subseteq P.T$ be the subsequence of $P.T$ which conforms to K, L, G with the smallest size. Note that there could be many qualified T' 's. Let the i^{th} local-consecutive part of T' be l_i and let the i^{th} gap of T' be g_i . Then, the size of T' can be written as $\sum_i (l_i + g_i)$. Since T' conforms to K, L, G , then $K + L \geq \sum_i (l_i) \geq K$, $l_i \geq L$, $g_i \leq G$. Therefore, $\sum_i (l_i + g_i) \leq \lfloor \frac{K}{L} \rfloor * G + K + L$. Thus ensuring each Par_t to be of that size would capture at least one of the T' 's, therefore the pattern P would valid in Par_t . This proves the completeness of the partitioning method. \square

5.0.2 Line Sweep Mining

After partition, each task in the reduce phase processes a partition Par_i , which contains $\lfloor \frac{K}{L} \rfloor * G + K + L$ snapshots starting from snapshot S_i . With such a partition method, we observe that within Par_i , only the patterns whose object sets are contained in the first snapshot are necessary to be reported. Therefore, we design a simple *line-sweep mining* (LSM) method to discover GCMPs. The algorithm works as in Algorithm 2.

The algorithm scans the snapshots in a partition in the sequential order. During the scan it maintains a candidate set

Algorithm 2 Line Sweep Mining

Require: $Par_t = \{S_t, S_{t+1}, \dots\}$

- 1: $C \leftarrow \{\}$ ▷ Candidate set
- 2: **for** $c \in S_t$ **do**
- 3: $C.add(\langle c, t \rangle)$
- 4: **end for**
- 5: **for** $i = 1; i < |Par_t|; i++$ **do**
- 6: $N \leftarrow S_i \oplus C$
- 7: **for all** $n \in N$ **do**
- 8: **if** $|n.O| \geq M$ **then** $C.add(n)$.
- 9: **end if**
- 10: **end for**
- 11: remove unqualified candidate from C
- 12: **end for**
- 13: output qualified candidate in C

C which could potentially be a valid pattern (line 1). The algorithm starts by inserting clusters at S_1 to C (lines 2-4). Subsequently, in each iteration, clusters in C are joined with clusters at S_i to generate a new set of patterns N (lines 6). Any valid new patterns are inserted back to C and any invalid patterns are discarded (lines 8 and 11).

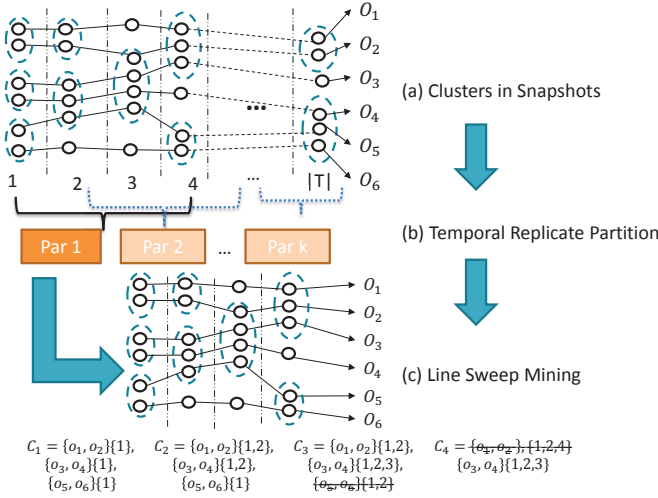


Figure 4: Work flow of trajectory replication and mining

Example 1. We illustrate the process of Trajectory Replication and Mining in Figure 4 using $M = 2, K = 2, L = 2, G = 2$. In (a), snapshots are clustered and these snapshots are the input to the TRM. Then, we compute the size for each partition, which equals to $\lfloor \frac{K}{L} - 1 \rfloor * G + K + L = 4$. Therefore, in (b), every four snapshots are grouped into a partition. Then a line sweep method is performed in (c) for partition 1. Each C_i refers to the candidate set when the algorithm sweeps each snapshot. Initially, C_1 contains patterns whose object set is in S_1 . When scanning the other snapshots, patterns in C_1 grow their timestamps. At S_3 , since the timestamps of $\{o_5, o_6\}$ is $\{1, 2\}$ which is neither a qualified pattern nor matches G constraint, thus this candidate is removed from C_3 . When line sweep finishes, only $\{o_3, o_4\}$ is the qualified pattern and is outputted.

The TRM approach though achieves good parallelism, it requires to replicate the data multiple times. Specifically,

each snapshots are copied $\lfloor \frac{K}{L} - 1 \rfloor * G + K + L$ times. In the cases of *swarm*, *group* and *platoon*, G is as large as $|T|$. Handling those cases is equivalent to replicate the entire snapshots to each partition, which surrenders the benefit of parallelism.

6. STAR PARTITION AND MINING

In order to develop a method that achieves good parallelism under any pattern parameters, we propose the *Star Partition and Mining* (SPM) method. In SPM, we first design a novel object-based partition method named *star partition*. Then, we design an *Apriori*-like method to mine GCMP patterns out of each partition. The overview of the SPM method is presented in Algorithm 3. As shown, the SPM method takes three phases. In the map phase, objects from the same cluster forms object-object pairs. The object-object pairs are then paired up with the timestamp of the snapshot to form a 3-uple (lines 1-8). In the partition phase, 3-uples with the same leading object form a *star* which will be explained shortly (lines 9-12). Lastly in the reduce phase, patterns are mined from each star structure (lines 13-16).

Algorithm 3 Star Partition and Mining

Require: list of $\langle t, S_t \rangle$ pairs

- 1: —Map phase—
- 2: **for all** $C \in S_t$ **do**
- 3: **for all** $(o_1, o_2) \in C \times C$ **do**
- 4: **if** $o_1 < o_2$ **then**
- 5: emit a $\langle o_1, o_2, \{t\} \rangle$ triplet
- 6: **end if**
- 7: **end for**
- 8: **end for**
- 9: —Partition and Shuffle phase—
- 10: **for all** $\langle o_1, o_2, \{t\} \rangle$ triplets **do**
- 11: group-by o_1 , emit $\langle o_1, Sr_{o_1} \rangle$
- 12: **end for**
- 13: —Reduce phase—
- 14: **for all** $\langle o, Sr_o \rangle$ **do**
- 15: Apriori(Sr_o)
- 16: **end for**

6.0.3 Star Partition

The intuition of star partition is that, if two objects are part of the same pattern, then at some snapshots, they must belong to the same cluster. Therefore, we may link objects that belong to the same cluster at some snapshots to form a *connection graph*. Objects that are not connected surely do not belong to a pattern. We may further partition the connection graph so that mining proper patterns can be done in parallel. We define the *connection graph* and *star* as follows:

Definition 5 (Connection Graph and Star). A *connection graph* is an undirected graph $G = (V : E)$, where each $v \in V$ represents an object. An edge $e(s, t) = ET \in E$ contains all the timestamps at which s, t are in the same cluster, i.e., $\forall t \in ET, C_t(s) = C_t(t)$. A *star* of a vertex s , denoted as Sr_s , is the set of incidental edges on s whose another ending vertex is greater than s . I.e., $\forall e(s, t) \in Sr_s, s < t$. We name s as the central vertex of Sr_s .

It is notable that we require vertices in a star to be greater than its central vertex. This effectively avoids replicating

edges. *Connection graph* and *star* examples are shown in Figure 5 (a) and (b). In (a), a connection graph is formed based on the example in Figure 1. In (b), 5 stars are presented. It is easy to see that, by requiring the center vertex to be the smallest vertex in a star, there are no edges been replicated. In implementation, as stated in Algorithm 3 line 4, the comparison between vertices/objects are based on the vertex/object ID.

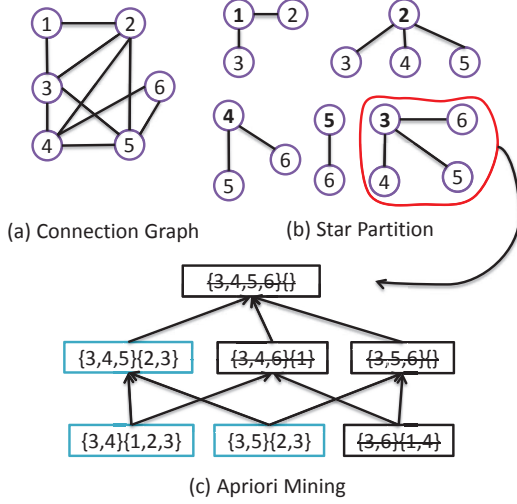


Figure 5: Star partition and mining of trajectories in Figure 1

Although star partition is performed based on the object connections, each star can be effectively viewed as a subset of trajectories. To see this, each vertex in a star can be viewed as an object. The timestamps of center vertex s is the union of all the edges in Sr_s . The timestamps of vertex $v \neq s$ is the edge $e(s, v)$. Therefore, we are able to define and mine GCMP on the stars. Before describing the mining strategy, we first state in the following theorem that the star-partition is complete and sound:

Theorem 2 (Soundness and Completeness of Star Partition). *Star partition is sound and complete.*

Proof. For the soundness, if P is a valid pattern in Sr_s , then at every time t , $\forall o_1, o_2 \in P.O$, $C_t(o_1) = C_t(o_2)$. By definition, P is valid in the original trajectories. For the completeness, if P is a valid pattern in original trajectories, let s be the object with smallest ID in $P.O$. Then by the definition of pattern, $\forall t \in P.T$, $\forall o \in P.O$, $C_t(s) = C_t(o)$. It follows that all object $o \in P.O$ are in Sr_s . Furthermore, every timestamp in $P.T$ is included in Sr_s . Therefore, P is a valid pattern in Sr_s . \square

Based on the above theorem, we can mine GCMP from each partition independently. It is notable that, in star partition, original data is replicated for $O(\mathbb{O})$ times as each object may be sent to $O(\mathbb{O})$ partitions. Since this complexity is free from pattern parameters, the star partition is more scalable than the temporal replication. In later sections, we will describe several engineering level optimization to further reduce the amount of replicated data.

6.0.4 Apriori Mining

In the mining phase, we need to find the patterns within each star. To systematically discover the patterns, we design the *Apriori Mining* method which is similar to the technique in frequent item mining literature. During the algorithm, we call a candidate pattern R -pattern if the size of its object set is R . Our algorithm runs in iterations. During each iteration R , we try to generate all $(R + 1)$ -patterns. In iteration 1, the 2-pattern is the edges in Sr_s . In particular, for each $e(s, v) = ET$, pattern $p = (\{s, v\}, ET)$ is formed. During each iteration, we generate $(R + 1)$ -patterns by joining R -patterns with 2-patterns. Specifically, the join between $p_1 = (O_1, T_1)$ and $p_2 = (O_2, T_2)$ would generate a new pattern $p_3 = (O_1 \cup O_2, T_1 \cap T_2)$. Notice that in Sr_s , each R -pattern consists of the object s , thus the join will grow a R -pattern at most to a $(R + 1)$ -pattern. Our mining algorithm stops where no further patterns are generated. The algorithm is illustrated as in Algorithm 4.

Algorithm 4 Apriori Mining

Require: Sr_s

```

1:  $Lv \leftarrow \{\}$ 
2:  $Ground \leftarrow \{\}$ 
3: for all  $e(s, t) = T \in Sr_s$  do
4:    $Ground.add(\{s, t\}, T)$ ;
5:  $Lv \leftarrow Ground$ ;
6: end for
7: while true do
8:   if  $Lv$  is not empty then
9:      $LvCand \leftarrow \{\}$ 
10:    for all  $cand_v \in Lv$  do
11:      for all  $cand \in Ground$  do
12:         $p \leftarrow cand_v \text{ join } cand$ 
13:        if  $p.T$  is partly valid then
14:           $LvCand.add(p)$ 
15:        end if
16:      end for
17:    end for
18:     $Lv \leftarrow LvCand$ 
19:  else
20:    break
21:  end if
22: end while
```

An illustration of Algorithm 4 is shown in Figure 5 (c). As shown, the star $Sr_3 = \{3, 4, 5, 6\}$ initially generate three 2-candidates. At every iteration, higher level candidates are generated by joining lower level candidates. When no more candidates can be generated, the algorithm stops by outputting the valid patterns.

It is notable that Algorithm 4 takes exponential complexity to mine GCMP. There are two major factors dragging down the performance. First, the size of Sr_s affects the initial size of 2-patterns. Second, the candidates generated in each level affects the join performance. In later sections, we exploit the property of GCMP to reduce the two factors.

7. OPTIMIZATION

In this section, we describe several optimizations to the star-partition and mining algorithm. In addition, we also address some practical issues when deploying the SPM algorithm to real MapReduce based systems.

7.1 Edge Simplification

Each edge $e(s, v)$ in Sr_s contains a time sequence ET which represents the co-occurrence of s and v . We notice that the edge between s and t is not always necessary. For example, if an edge has a cardinality less than K , it is unnecessary to include this edge to Sr_s since it cannot contribute to any patterns. This motivates us to simplify the edges in Sr_s to boost the overall performance.

Our goal of edge simplification is to, given a time sequence T , find a subsequence of $T' \subseteq T$, such that T' is potentially conforms to K, L, G . And we wish $|T'|$ to be as small as possible. We start-off by observing that for every time sequence T , T can be divided into a set of maximally G -connected subsequences. Note that a maximally G -connected subsequence can potentially contribute to a pattern if it conforms to K, L . Therefore, we are able to reduce T to its maximally G -connected subsequences which conform to K, L .

To formally describe the idea, we define the a *candidate sequence* as follows:

Definition 6 (Candidate Sequence). *Given the pattern parameters: L, K, G , a sequence T is a Candidate Sequence if for any of its maximal G -connected sequence T' , T' conforms to L, K .*

For example, let $L = 2, K = 4, G = 2$, sequence $T_1 = (1, 2, 4, 5, 6, 9, 10, 11, 13)$ is not a fully candidate sequence since one of its maximal G -connected sequence $(9, 10, 11)$ is not a partly candidate sequence. In contrast, sequence $T_2 = (1, 2, 4, 5, 6)$ is a fully candidate sequence.

To reduce a sequence T to a candidate sequence, we need to strip out its maximal G -connected subsequences which does not form to K, L . Such a reduction takes two rounds scan of T as shown in Algorithm 5. In the first round, the consecutive portions of T with size less than L are removed. In the second round, the maximal G -connected sequences of size less than K are removed. Clearly the simplification algorithm runs in $O(|T|)$ time.

Example 2. Take $T_1 = \{1, 2, 4, 5, 6, 9, 10, 11, 13\}$ as an example of edge simplification. Let $L = 2, K = 4, G = 2$. In the first round of scan. T_1 reduces to $\{1, 2, 4, 5, 6, 9, 10, 11\}$. The consecutive subsequence $\{13\}$ is removed by $L = 2$. T_1 has two maximal G -consecutive subsequences, which are $\{1, 2, 4, 5, 6\}$ and $\{9, 10, 11\}$. Since $K = 4$, $\{9, 10, 11\}$ is removed from T_1 in the second round of scan. Therefore, T_1 is simplified to $\{1, 2, 4, 5, 6\}$.

By leveraging the edge simplification technique, the size of the edges in Sr_s can be greatly reduced. If an edge cannot be reduced to a candidate sequence, then it is directly removed from Sr_s . If an edge can be reduced to a candidate sequence, replacing itself by the candidate sequence results in a more compact storage.

7.2 Candidate Pruning via Temporal Monotonicity

During the apriori phase, we repeatedly join candidate patterns in different levels to generate a larger set of a patterns. We observe that traditional monotonic property of Apriori algorithms **does not** hold in GCMP mining. That is given two candidate P_1, P_2 , if $P_1.O \subset P_2.O$ and P_1 is not a valid pattern, then P_2 may or may not be a valid pattern. However, we notice that we may form another monotonic property based on the *candidate sequence* such that the Apriori algorithm could still benefit.

Algorithm 5 Edge Simplification

Require: T

- 1: —Remove the consecutive portion with size less than L —
- 2: $c \leftarrow 0$
- 3: **for** $i \in (0, \dots, |T|)$ **do**
- 4: **if** $T[i] - T[i - 1] = 1$ **then**
- 5: **if** $i - c < L$ **then**
- 6: T remove $[c : i)$
- 7: **end if**
- 8: $c \leftarrow i$
- 9: **end if**
- 10: **end for**
- 11: —Remove the pseudo-consecutive portion with size less than K —
- 12: $s \leftarrow 1, c \leftarrow 0$
- 13: **for** $i \in (0 : |T|)$ **do**
- 14: **if** $T[i] - T[i - 1] > G$ **then**
- 15: **if** $s < K$ **then**
- 16: T remove $[c : i)$
- 17: **end if**
- 18: $c \leftarrow i, s \leftarrow 1$
- 19: **else**
- 20: $s++$
- 21: **end if**
- 22: **end for**

The intuition is that if a candidate $P_1.T$ cannot be reduced to a *candidate sequence*, then P_1 cannot be valid pattern. Furthermore, any candidate P_2 , with $P_1.O \subset P_2.O$ cannot be a valid pattern. This *temporal monotonic property* is explicitly described as in the follow theorem:

Theorem 3 (Temporal Monotonic Property of GCMP). *Given the temporal parameters L, G, K , for a candidate c in Algorithm 4, if $c.T$ cannot be reduced to a candidate sequence, then for any candidate c' with $c.O \subset c'.O$, c' can be pruned.*

Proof. Let c_1, c_2 be two candidates with $c_1.O \subset c_2.O$. It is easy to see that $c_1.T \supseteq c_2.T$. If $c_1.T$ cannot be reduced to a candidate sequence, then any subset of $c_1.T$ cannot be reduced. It follows that $c_2.T$ cannot be reduced neither. Thus, if $c_1.T$ cannot be reduced to a candidate sequence, c_2 can be pruned. \square

Example 3. We use Figure 5 (c) to demonstrate the candidate pruning. As shown, at the initial stage, $\{3, 6\}\{1, 4\}$ is pruned, since $\{1, 4\}$ is not a candidate sequence. By temporal monotonicity, candidates containing objects $\{3, 6\}$ can all be pruned. Therefore, we are able to directly prune $\{3, 4, 6\}$, $\{3, 5, 6\}$ and $\{3, 4, 5, 6\}$.

7.3 Load Balancing

7.4 Duplication Detection

7.5 Handling Overlapping Clusters

8. EXPERIMENTAL STUDY

9. CONCLUSION AND FUTURE WORK

10. REFERENCES

- [1] Y. Zheng, Y. Liu, J. Yuan, and X. Xie, “Urban computing with taxicabs,” in *Proceedings of the 13th international conference on Ubiquitous computing*, pp. 89–98, ACM, 2011.
- [2] Z. Li, B. Ding, J. Han, R. Kays, and P. Nye, “Mining periodic behaviors for moving objects,” in *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 1099–1108, ACM, 2010.
- [3] J. Bao, Y. Zheng, D. Wilkie, and M. F. Mokbel, “A survey on recommendations in locationbased social networks. submitted to,” *Geoinformatica*, 2013.
- [4] X. Li, *Managing moving objects and their trajectories*. PhD thesis, National University of Singapore, 2013.
- [5] J. Gudmundsson and M. van Kreveld, “Computing longest duration flocks in trajectory data,” in *Proceedings of the 14th annual ACM international symposium on Advances in geographic information systems*, pp. 35–42, ACM, 2006.
- [6] Y. Wang, E.-P. Lim, and S.-Y. Hwang, “Efficient mining of group patterns from user movement data,” *Data & Knowledge Engineering*, vol. 57, no. 3, pp. 240–282, 2006.
- [7] H. Jeung, M. L. Yiu, X. Zhou, C. S. Jensen, and H. T. Shen, “Discovery of convoys in trajectory databases,” *Proceedings of the VLDB Endowment*, vol. 1, no. 1, pp. 1068–1080, 2008.
- [8] Z. Li, B. Ding, J. Han, and R. Kays, “Swarm: Mining relaxed temporal moving object clusters,” *Proceedings of the VLDB Endowment*, vol. 3, no. 1-2, pp. 723–734, 2010.
- [9] Y. Li, J. Bailey, and L. Kulik, “Efficient mining of platoon patterns in trajectory databases,” *Data & Knowledge Engineering*, 2015.
- [10] J. Gudmundsson, M. van Kreveld, and B. Speckmann, “Efficient detection of motion patterns in spatio-temporal data sets,” in *Proceedings of the 12th annual ACM international workshop on Geographic information systems*, pp. 250–257, ACM, 2004.
- [11] Y. Zheng, “Trajectory data mining: an overview,” *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 6, no. 3, p. 29, 2015.
- [12] D. Birant and A. Kut, “St-dbscan: An algorithm for clustering spatial-temporal data,” *Data & Knowledge Engineering*, vol. 60, no. 1, pp. 208–221, 2007.
- [13] P. Laube, M. van Kreveld, and S. Imfeld, “Finding remodetecting relative motion patterns in geospatial lifelines,” in *Developments in spatial data handling*, pp. 201–215, Springer, 2005.
- [14] D. H. Douglas and T. K. Peucker, “Algorithms for the reduction of the number of points required to represent a digitized line or its caricature,” *Cartographica: The International Journal for Geographic Information and Geovisualization*, vol. 10, no. 2, pp. 112–122, 1973. doi:10.3138/FM57-6770-U75U-7727.
- [15] J. Dean and S. Ghemawat, “Mapreduce: simplified data processing on large clusters,” *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.